

**AMENDMENTS TO THE CLAIMS:**

This listing of claims will replace all prior versions and listings of claims in the application:

**Listing of Claims:**

1. **(Currently amended)** A method of clustering documents [[I]] or patterns [[D]] each having one or plural document [[I]] or pattern [[D]] segments in an input document [[I]] or pattern [[D]] set, said method comprising: ~~based on a relation among them, comprising,~~

(a) obtaining a document [[I]] or pattern [[D]] frequency matrix for the set of input documents [[I]] or patterns [[D]], based on occurrence frequencies of terms appearing in each document [[I]] or pattern [[D]];

(b) selecting a seed document [[I]] or pattern [[D]] from remaining documents [[I]] or patterns [[D]] that are not included in any cluster existing at that moment and constructing a current cluster of the initial state using the seed document [[I]] or pattern [[D]];

(c) obtaining the document [[I]] or pattern [[D]] commonality to the current cluster for each document [[I]] or pattern [[D]] in the input document [[I]] or pattern [[D]] set by using information based on the document [[I]] or pattern [[D]] frequency matrix for the input document [[I]] or pattern [[D]] set, information based on the document [[I]] or pattern [[D]] frequency matrix for documents [[I]] or patterns [[D]] in the current cluster and information based on the common co-occurrence matrix of the current cluster, and making documents [[I]] or patterns [[D]] having the document commonality higher than a threshold belong temporarily to the current cluster;

(d) repeating step (c) until the number of documents [[I]] or patterns [[D]] temporarily belonging to the current cluster becomes the same as that in the previous repetition;

(e) repeating steps (b) through (d) until a given convergence condition is satisfied; and

(f) deciding, on the basis of the document [[I]] or pattern [[D]] commonality of each document [[I]] or pattern [[D]] to each cluster, a cluster to which each document [[I]] or pattern [[D]] belongs and outputting said cluster.

2. **(Currently amended)** A clustering method according to claim 1, wherein step (a) further includes; [[.]]

(a-1) generating a document [[I]] or pattern [[D]] segment vector for each of said document [[I]] or pattern [[D]] segments based on occurrence frequencies of terms appearing in each document [[I]] or pattern [[D]] segment;

(a-2) obtaining a co-occurrence matrix for each document [[I]] or pattern [[D]] in the input document [[I]] or pattern [[D]] set from the document [[I]] or pattern [[D]] segment vectors; and

(a-3) obtaining a document [[I]] or pattern [[D]] frequency matrix from the co-occurrence matrix for each document.

3. **(Currently amended)** A clustering method according to claim 1, wherein step (b) further includes; [[.]]

(b-1) constructing a common co-occurrence matrix of remaining documents [[I]] or patterns [[D]] that are not included in any cluster existing at that moment; and

(b-2) obtaining a document commonality to the set of the remaining document [(I)] or pattern [(D)] set for each document [(I)] or pattern [(D)] in the remaining document [(I)] or pattern [(D)] set by using the common co-occurrence matrix of the remaining documents [(I)] or patterns [(D)], and extracting the document [(I)] or pattern [(D)] having the highest document [(I)] or pattern [(D)] commonality, and constructing a current cluster of the initial state by making a document [(I)] or pattern [(D)] set including the seed document [(I)] or pattern [(D)] and the neighbor documents [(I)] or patterns [(D)] similar to the seed document [(I)] or pattern [(D)].

4. **(Currently amended)** A clustering method according to claim 1, wherein step (c) further includes: [(,)]

(c-1) constructing a common co-occurrence matrix of the current cluster and a document [(I)] or pattern [(D)] frequency matrix of the current cluster;

(c-2) obtaining the distinctiveness of each term and each term pair to the current cluster by comparing the document [(I)] or pattern [(D)] frequency matrix of the input document [(I)] or pattern [(D)] set and the document [(I)] or pattern [(D)] frequency matrix of the current cluster; and

(c-3) obtaining document [(I)] or pattern [(D)] commonalities to the current cluster for each document [(I)] or pattern [(D)] in the input document [(I)] or pattern [(D)] set by using the common co-occurrence matrix of the current cluster and weights of each term and term pair obtained from their distinctiveness, and making a document [(I)] or pattern [(D)] having the document [(I)] or pattern [(D)] commonality higher than a threshold belong temporarily to the current cluster.

5. **(Currently amended)** A clustering method according to claim 1, further including: [[.]]

repeating step (e) until the number of documents [[(I)] or patterns [(D)]] whose document [[(I)] or pattern [(D)]] commonalities to any current clusters are less than a threshold becomes 0, or the number is less than a threshold and is equal to that of the previous repetition.

6. **(Currently amended)** A clustering method according to claim 1, wherein step (f) further includes: [[.]]

checking existence of a redundant cluster, and removing, when the redundant cluster exists, the redundant cluster and again deciding the cluster to which each document belongs.

7. **(Currently amended)** A method according to claim 1, wherein the co-occurrence matrix  $S^r$  of the document [[(I)] or pattern [(D)]]  $D_r$  is determined in accordance with:

$$S^r = \sum_{y=1}^{Y_r} d_{ry} d_{ry}^T \quad (1)$$

where: M equals the number of sorts of the occurring terms,  $D_r$  equals the  $r$ th document [[(I)] or pattern [(D)]] in a document [[(I)] or pattern [(D)]] set D consisting of R documents [[(I)] or patterns [(D)]],  $Y_r$  equals the number of document [[(I)] or pattern [(D)]] segments in document [[(I)] or pattern [(D)]]  $D_r$ , and  $d_{ry} = [[(I)] d_{ry1}, \dots, d_{ryM} [(D)]]^T$  equals the  $y$ th document [[(I)] or pattern [(D)]] segment vector of document [[(I)] or pattern [(D)]]  $D_r$ , and T represents transposition of a vector.

8. **(Currently amended)** A method according to claim 1, wherein each component of the document [[(I)] or pattern [(D)]] frequency matrix of a document [[(I)] or pattern [(D)]] set D is

the number of documents  $[[I]]$  or patterns  $[[D]]$  in which a corresponding component of the co-occurrence matrix of each document  $[[I]]$  or pattern  $[[D]]$  in the document  $[[I]]$  or pattern  $[[D]]$  set D does not take a value of zero.

9. **(Currently amended)** A method according to claim 1, further comprising:  
determining the common co-occurrence matrix of a document  $[[I]]$  or pattern  $[[D]]$  set D from a matrix  $T^A$  on the basis of a matrix T whose mn component is determined by the matrix  $T^A$  having an mn component determined by

$$\begin{aligned} T_{mn}^A &= T_{mn}, & U_{mn} > A, \\ T_{mn}^A &= 0 & \text{otherwise,} \end{aligned}$$

where  $U_{mn}$  represents the mn component of the document  $[[I]]$  or pattern  $[[D]]$  frequency matrix of the document  $[[I]]$  or pattern  $[[D]]$  set D.

10. **(Currently amended)** A method according to claim 1, further comprising:  
determining the common co-occurrence matrix of a document  $[[I]]$  or pattern  $[[D]]$  set D from a matrix  $Q^A$  on the basis of a matrix T whose mn component is determined by

$$\begin{aligned} T_{mn} &= \prod_{r=1}^R S_{mn}^r \\ S_{mn}^r &> 0 \end{aligned}$$

the matrix  $Q^A$  having an mn component determined by

$$\begin{aligned} Q_{mn}^A &= \log[[I]] T_{mn}^A [[D]] & T_{mn}^A > 1, \\ Q_{mn}^A &= 0 & \text{otherwise.} \end{aligned}$$

11. **(Currently amended)** A method according to claim 10<sub>a</sub> wherein  $z_{mm}$  and  $z_{mn}$  are respectively weights for a term  $[(I)]$  or object feature  $[(D)]$   $m$  and a term  $[(I)]$  or object feature  $[(D)]$  pair  $m, n$ , a document  $[(I)]$  or pattern  $[(D)]$  commonality of document  $[(I)]$  or pattern  $[(D)]$   $P$  having a co-occurrence matrix  $S^P$  with respect to the document  $[(I)]$  or pattern  $[(D)]$  set  $D$  given by

$$com_l(D, P; Q^A) = \frac{\sum_{m=1}^M z_{mm} Q^A_{mm} S^P_{mm}}{\sqrt{\sum_{m=1}^M z_{mm} (Q^A_{mm})^2} \sqrt{\sum_{m=1}^M z_{mm} (S^P_{mm})^2}} \quad (3)$$

or

$$com_q(D, P; Q^A) = \frac{\sum_{m=1}^M \sum_{n=1}^M z_{mn} Q^A_{mn} S^P_{mn}}{\sqrt{\sum_{m=1}^M \sum_{n=1}^M z_{mn} (Q^A_{mn})^2} \sqrt{\sum_{m=1}^M \sum_{n=1}^M z_{mn} (S^P_{mn})^2}} \quad (4).$$

12. **(Currently amended)** A method according to claim 9<sub>a</sub> wherein  $z_{mm}$  and  $z_{mn}$  are respectively weights for a term  $[(I)]$  or object feature  $[(D)]$   $m$  and a term  $[(I)]$  or object feature  $[(D)]$  pair  $m, n$ , a document  $[(I)]$  or pattern  $[(D)]$  commonality of document  $[(I)]$  or pattern  $[(D)]$   $P$  having a co-occurrence matrix  $S^P$  with respect to the document  $[(I)]$  or pattern  $[(D)]$  set  $D$  given by

$$com_l(D, P; T^A) = \frac{\sum_{m=1}^M z_{mm} T^A_{mm} S^P_{mm}}{\sqrt{\sum_{m=1}^M z_{mm} (T^A_{mm})^2} \sqrt{\sum_{m=1}^M z_{mm} (S^P_{mm})^2}} \quad (3)$$

or

$$com_q(D, P; T^A) = \frac{\sum_{m=1}^M \sum_{n=1}^M z_{mn} T^A_{mn} S^P_{mn}}{\sqrt{\sum_{m=1}^M \sum_{n=1}^M z_{mn} (T^A_{mn})^2} \sqrt{\sum_{m=1}^M \sum_{n=1}^M z_{mn} (S^P_{mn})^2}} \quad (4).$$

13. **(Currently amended)** A method according to claim 1, wherein extraction of the seed document [(I)] or pattern [(D)] of the current cluster and construction of the current cluster of the initial state ~~includes~~ comprise:

(a') obtaining a document [(I)] or pattern [(D)] commonality to the remaining document [(I)] or pattern [(D)] set for each document [(I)] or pattern [(D)] in the remaining document [(I)] or pattern [(D)] set by using the said common co-occurrence matrix of the remaining documents [(I)] or patterns [(D)],

(b') extracting, as candidates of the seed of the current cluster, a specific number of documents [(I)] or patterns [(D)] whose document [(I)] or pattern [(D)] commonalities obtained by step (a') are large;

(c') obtaining similarities of the respective candidates of the seed of the cluster to all documents [(I)] or patterns [(D)] in the input document [(I)] or pattern [(D)] set or in the remaining document [(I)] or pattern [(D)] set, and obtaining documents [(I)] or patterns [(D)] having similarities larger than a threshold as neighbor documents [(I)] or patterns [(D)] of the candidate; and

(d') selecting the candidate whose number of the neighbor documents [(I)] or patterns [(D)] is the largest among the candidates as the seed of the current cluster and making its neighbor documents [(I)] or patterns [(D)] the current cluster of the initial state.

14. **(Currently amended)** A method according to claim 1, further including:

detecting the distinctiveness of each term [(I)] or object feature [(D)] and each term pair with respect to the current cluster and detecting their weights,

the distinctiveness and weight detecting steps including:

(a'') obtaining a ratio of each component of a document  $[[I]]$  or pattern  $[[D]]$  frequency matrix obtained from the input document  $[[I]]$  or pattern  $[[D]]$  set to a corresponding component of a document  $[[I]]$  or pattern  $[[D]]$  frequency matrix obtained from the current cluster as a document  $[[I]]$  or pattern  $[[D]]$  frequency ratio of each term  $[[I]]$  or feature  $[[D]]$  or each term  $[[I]]$  or feature  $[[D]]$  pair;

(b'') selecting a specific number of terms  $[[I]]$  or features  $[[D]]$  or term  $[[I]]$  or feature  $[[D]]$  pairs having the smallest document  $[[I]]$  or pattern  $[[D]]$  frequency ratios among a specific number of terms  $[[I]]$  or features  $[[D]]$  or term  $[[I]]$  or feature  $[[D]]$  pairs having the highest document  $[[I]]$  or pattern  $[[D]]$  frequencies, and obtaining the average of the document  $[[I]]$  or pattern  $[[D]]$  frequency ratios of the selected terms  $[[I]]$  or features  $[[D]]$  or term  $[[I]]$  or feature  $[[D]]$  pairs as the average document  $[[I]]$  or pattern  $[[D]]$  frequency ratio;

(c'') dividing the average document  $[[I]]$  or pattern  $[[D]]$  frequency ratio by the document  $[[I]]$  or pattern  $[[D]]$  frequency ratio of each term  $[[I]]$  or feature  $[[D]]$  or each term  $[[I]]$  or feature  $[[D]]$  pair as a measure of the distinctiveness of each term  $[[I]]$  or feature  $[[D]]$  or each term  $[[I]]$  or feature  $[[D]]$  pair; and

(d'') determining the weight of each term  $[[I]]$  or feature  $[[D]]$  or each term  $[[I]]$  or feature  $[[D]]$  pair from a function having the distinctiveness measure as a variable.

15. **(Currently amended)** A method according to claim 1<sub>a</sub> further including:  
eliminating terms  $[[I]]$  or features  $[[D]]$  and term  $[[I]]$  or feature  $[[D]]$  pairs having document  $[[I]]$  or pattern  $[[D]]$  frequencies higher than a threshold.



16. **(Currently amended)** A method according to claim 1, wherein clustering is performed recursively by letting the document  $[[I]]$  or pattern  $[[D]]$  set included in a cluster be the input document  $[[I]]$  or pattern  $[[D]]$  set.

17. **(Currently amended)** A computer program product containing a computer program which, when executed by a computer, causes the ~~for causing a~~ computer to perform the method of claim 1.

18. **(Currently amended)** A computer program product containing a computer program which, when executed by a computer, causes the ~~for causing a~~ computer to perform the method of claim 2.

19. **(Currently amended)** A computer program product containing a computer program which, when executed by a computer, causes the ~~for causing a~~ computer to perform the method of claim 3.

20. **(Currently amended)** A computer program product containing a computer program which, when executed by a computer, causes the ~~for causing a~~ computer to perform the method of claim 4.

21. **(Currently amended)** A computer program product containing a computer program which, when executed by a computer, causes the ~~for causing a~~ computer to perform the method of claim 5.

22. **(Currently amended)** A computer program product containing a computer program which, when executed by a computer, causes the ~~for causing a~~ computer to perform the method of claim 6.

23. (Original) A computer arranged to perform the method of claim 1.

24. (Original) A computer arranged to perform the method of claim 2.

25. (Original) A computer arranged to perform the method of claim 3.

26. (Original) A computer arranged to perform the method of claim 4.

27. (Original) A computer arranged to perform the method of claim 5.

28. (Original) A computer arranged to perform the method of claim 6.

29. **(Currently amended)** A clustering apparatus for clustering documents [(I)] or patterns [(D)] each having one or plural document [(I)] or pattern [(D)] segments in an input document [(I)] or pattern [(D)] set ~~based on the relation among them~~, the apparatus comprising:

[(a)] a first unit [(means)] for obtaining a document [(I)] or pattern [(D)] frequency matrix for the set of input documents [(I)] or patterns [(D)], based on occurrence frequencies of terms appearing in each document [(I)] or pattern [(D)];

[(b)] a second unit [(means)] for selecting a seed document [(I)] or pattern [(D)] from remaining documents [(I)] or patterns [(D)] that are not included in any cluster existing at that moment and constructing a current cluster of the initial state using the seed document [(I)] or pattern [(D)];

[(c)] a third unit [(means)]

for obtaining the document [(I)] or pattern [(D)] commonality to the current cluster for each document [(I)] or pattern [(D)] in the input document [(I)] or pattern [(D)] set using information based on the document [(I)] or pattern [(D)] frequency matrix for the input document [(I)] or pattern [(D)] set, information based on the document [(I)] or pattern [(D)] frequency matrix for documents [(I)] or patterns [(D)] in the current cluster and information based on the common co-occurrence matrix of the current cluster, and [(means)]

for making documents [(I)] or patterns [(D)] having the document [(I)] or pattern [(D)] commonality higher than a threshold belong temporarily to the current cluster;

[(d)] a fourth unit [[means]] for repeating the operations of the third unit ~~means (e)~~ until the number of documents [(I)] or patterns [(D)] temporarily belonging to the current cluster becomes the same as that in the previous repetition;

[(e)] a fifth unit [[means]] for repeating the operations of the second through fourth units ~~means (b) through (d)~~ until given convergence conditions are satisfied; and

[(f)] a sixth unit [[means]] for deciding, on the basis of the document [(I)] or pattern [(D)] commonality of each document [(I)] or pattern [(D)] to each cluster, a cluster to which each document [(I)] or pattern [(I)] belongs, and for outputting said cluster.